

# The API Privacy Paradox

How LLM APIs Unintentionally Reveal Model Details

Sean Ren | Associate Professor, USC

✉ xiangren@usc.edu | ✉ @xiangrennlp

LLMs are often deployed behind API endpoints to keep the model's architecture and weights secret.

The assumption is that limited API access protects proprietary details... But is this true?

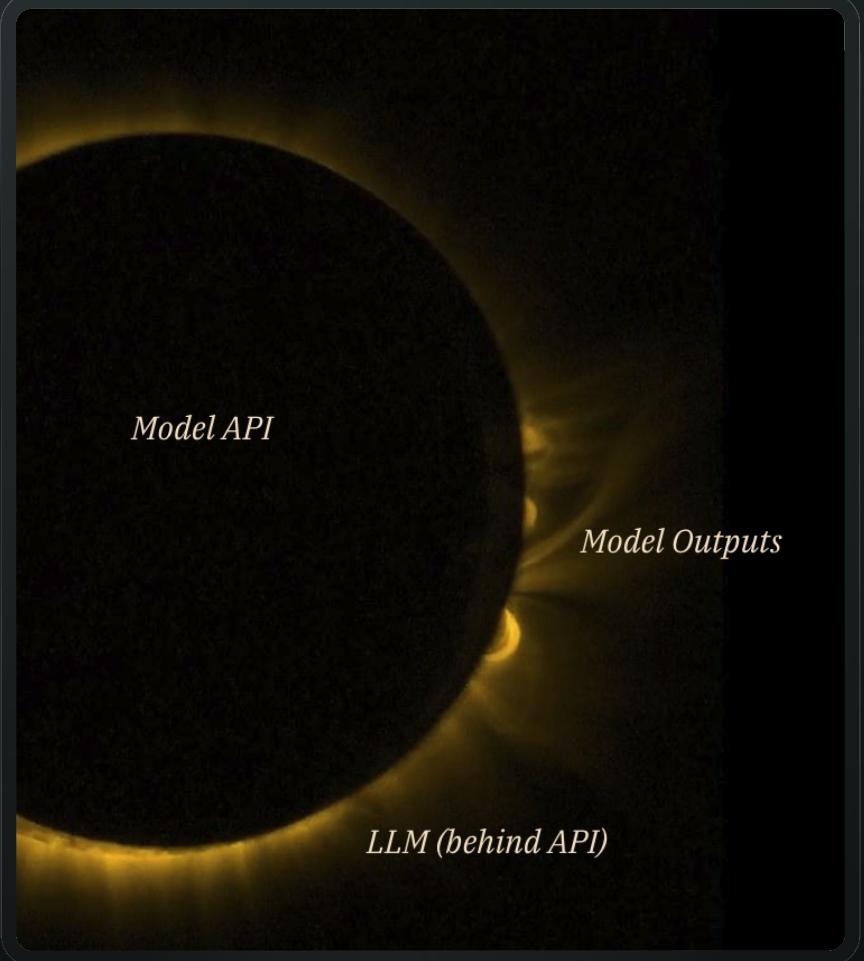
# System outputs reflect underlying structure

During an eclipse, scientists observe the structure of the corona (solar atmosphere) to study the magnetic topography of the temporarily hidden sun.



# What do LLM APIs Reveal?

- APIs don't just return text outputs; they can also return probabilities/logits.
- Probability distributions embed subtle architectural signals
- The Result: False sense of security as outputs can reveal hidden details



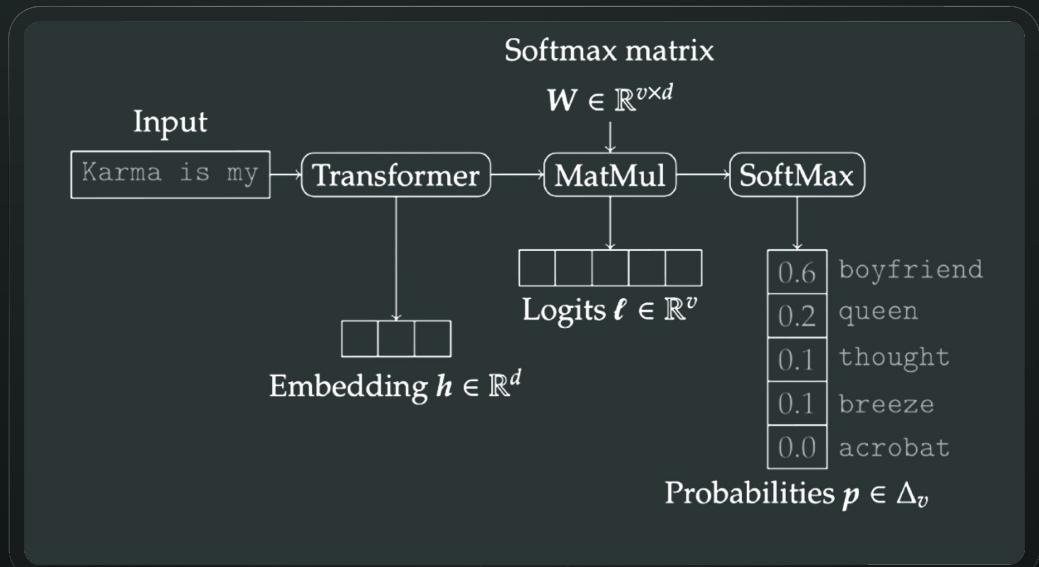
*Model API*

*Model Outputs*

*LLM (behind API)*

# Language Models

- When you give an LLM a prompt, it assigns a score (called a logit) to every possible next word
- These scores come from a fixed-size hidden state (embedding  $h$ )



# Language Models *Secrets*

## 2 Scope and Limitations of this Technical Report

This report focuses on the capabilities, limitations, and safety properties of GPT-4. GPT-4 is a Transformer-style model [39] pre-trained to predict the next token in a document, using both publicly available data (such as internet data) and data licensed from third-party providers. The model was then fine-tuned using Reinforcement Learning from Human Feedback (RLHF) [40]. Given both the competitive landscape and the safety implications of large-scale models like GPT-4, this report contains no further details about the architecture (including model size), hardware, training compute, dataset construction, training method, or similar.

Companies like OpenAI hide model details to protect their business.

**Secrets** include model parameters like the LM head and hyperparameters like the embedding size.

# API Assumptions

The mainstream APIs give top-k logprobs where k is usually 5:

TOP-K API OUTPUT

```
{  
    "isi": -2.73,  
    " some": -2.98,  
    ",": -3.03,  
    "1": -3.25,  
    " what": -3.43  
}
```

Logit bias lets the user “boost” the probability of specific tokens to put them in the top-k:

LOGIT BIAS

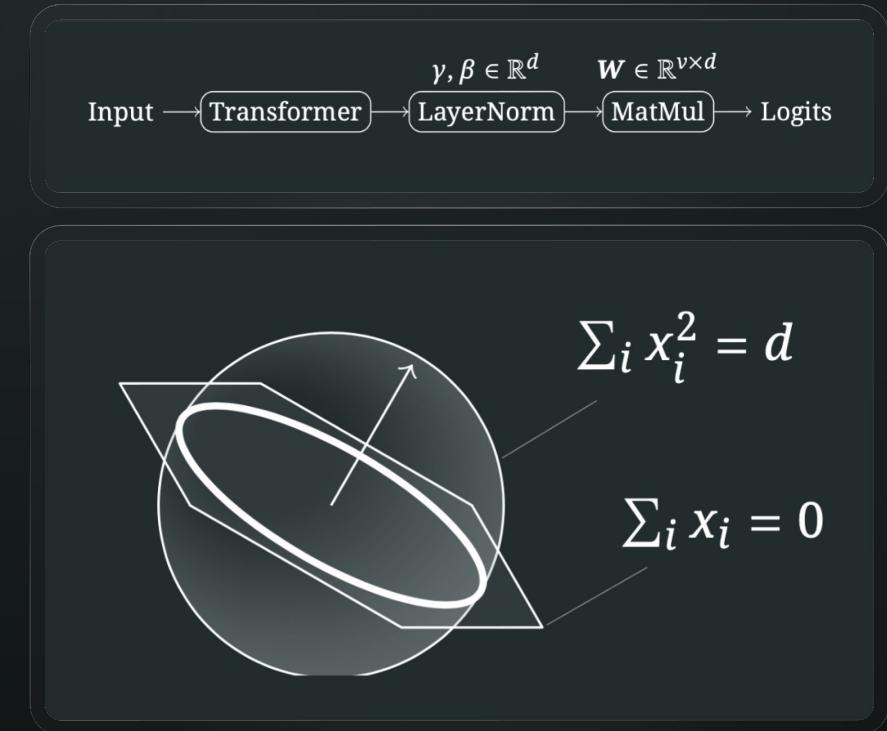
```
{" hi": +100}
```

BIASED TOP-K OUTPUT

```
{  
    " hi": -1.02,  
    "isi": -3.23,  
    ...  
}
```

# Why Logits Leak Info: The Softmax Bottleneck

- Because of that fixed size, the outputs can't cover the entire vocabulary patterns of the LLM (GPT-3.5 has a vocabulary of >50,000 tokens)
- They're confined to this smaller  $d$ -dimensional subspace  
→ the **softmax bottleneck**



# Exploiting the Bottleneck

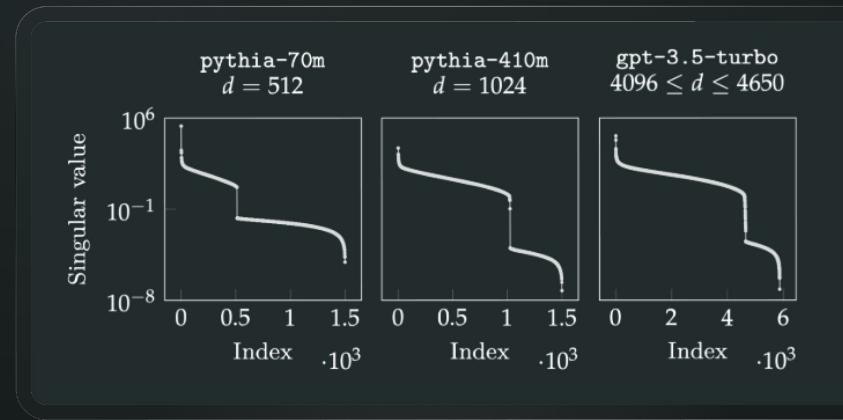
- Attackers send many diverse prompts to the API
- Collect the output probabilities/logits for each token
- Use SVD to find how many unique “directions” the outputs span
- The number of directions  $\approx$  the hidden size of the model

Cost estimation on GPT-3.5-Turbo  
API

Algorithm	Complexity	API calls per output	Image price (USD)
Logprob-free (Morris et al., 2023)	$v \log(\beta_{\max}/\epsilon)$	800 000	16 384
With logprobs	$v/k$	20 000	410

# Case Study: GPT3.5-turbo

- ~6,000 API queries
- Analysis showed outputs stopped adding new info after ~4,100 dimensions
- Hidden size inferred: ~4,096
- Enabled estimates of GPT3.5-Turbo's true scale without access to its weights



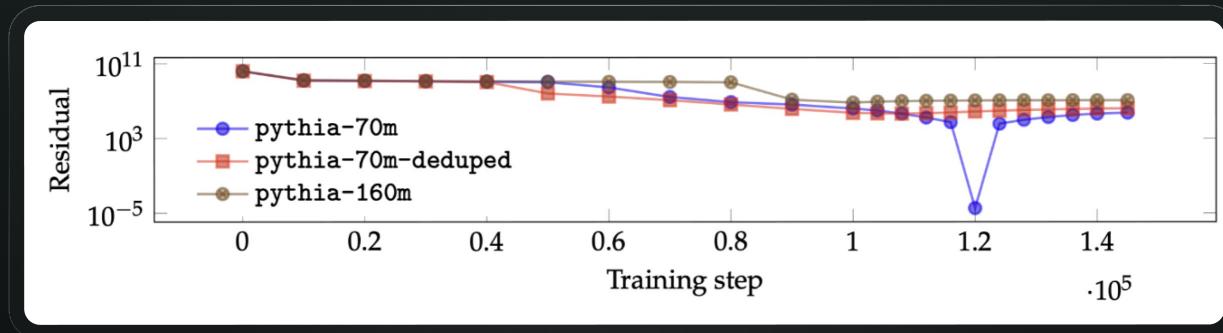
Applying the inverse log-softmax to the logprobs we get the logits  $L$ . All that remains is to check the rank of the logits  $(L)$ !

API outputs can reveal a model's true scale, letting outsiders estimate cost, capability, and even catch providers swapping or lying about their models.

# LLM images: threat or benefit?

APIs leak more than just the hidden size, they leak the column space of  $W$ . This can be useful for

- Identifying which model generated an output



# LLM images: threat or benefit?

APIs leak more than just the hidden size, they leak the column space of  $W$ . This can be useful for

- Auditing changes in an API LLM

Change	Interpretation
No logit change, no image change	No update
Logit change, no image change	Hidden prompt change or partial finetune
Low-rank image change (§7.1)	LoRA update
Image change	Full finetune

Table 2: Implications of image/logit changes.

# Security Risks from Softmax Leakage

## Risk 1

Enables model stealing via distillation or extraction attacks

## Risk 2

Lets competitors infer model scale and architecture

## Risk 3

Opens the door to crafting adversarial attacks using leaked behavior

## Risk 4

Low-cost feasibility as attacks can succeed for hundreds to a few thousand USD via API access

# Opportunities: Turning Leakage into Accountability

## Independent Audits

verify a provider's model size, safety claims, or fine-tune dates

## Change Detection

spot silent model swaps or unannounced updates in production

## Provenance Checks

confirm an API isn't a re-branded open-source model

## Transparency Signals

leakage can act as a public "fingerprint," building user trust

## Policy & Compliance

regulators can verify that licensed or safety-critical models stay unchanged

# Key Takeaways

- What do LLM APIs reveal about the model?
- What useful features can APIs keep without revealing model information?
- How can we verify the source of LLM outputs?

✉ xiangren@usc.edu | ✎ @xiangrennlp

## Logits of API-Protected LLMs Leak Proprietary Information

Matthew Finlayson Xiang Ren Swabha Swamydipta  
Thomas Lord Department of Computer Science  
University of Southern California  
[{mfinlays, xiangren, swabhas}@usc.edu](mailto:{mfinlays, xiangren, swabhas}@usc.edu)

### Abstract

Large language model (LLM) providers often hide the architectural details and parameters of their proprietary models by restricting public access to a limited API. In this work we show that, with only a conservative assumption about the model architecture, it is possible to learn a surprisingly large amount of non-public information about an API-protected LLM from a relatively small number of API queries (e.g., costing under \$1000 USD for OpenAI's gpt-3.5-turbo). Our findings are centered on one key observation: most modern LLMs suffer from a softmax bottleneck, which restricts the model outputs to a linear subspace of the full output space. We exploit this fact to unlock several capabilities, including (but not limited to) obtaining cheap full-vocabulary outputs, auditing for specific types of model updates, identifying the source LLM given a single full LLM output, and even efficiently discovering the LLM's hidden size. Our empirical investigations show the effectiveness of our methods, which allow us to estimate the embedding size of OpenAI's gpt-3.5-turbo to be about 4096. Lastly, we discuss ways that LLM providers can guard against these attacks, as well as how these capabilities can be viewed as a feature (rather than a bug) by allowing for greater transparency and accountability.



<https://mattf1n.github.io/>

# Future

How can we trust LLM providers?

How can we verify the source of LLM outputs?

How can we trace the contributions in AI development?

How to achieve fair compensation in a new AI economy?

# The Convergence of AI and Web3

Exploring the Infra Needed for a Collaborative, User-Owned AI Economy

# AI Is Already One of the Most Transformative Technologies Since the Internet...

- AI: **\$15.7T** projected impact on global GDP by 2030
- **78%** of enterprises now report using AI in at least one core function
- ChatGPT has **~800 million** weekly active users

# The Convergence of AI x Web3

- AI introduces **new digital primitives**: agents, datasets, and models that hold real-world value
- Web3 offers the tools for transparent ownership, on-chain attribution, and trustless collaboration

# The Convergence of AI x Web3

- AI introduces **new digital primitives**: agents, datasets, and models that hold real-world value
- Web3 offers the tools for transparent ownership, composability, and trustless collaboration
- Together, they unlock the possibility of building autonomous systems that are not only intelligent but governed and monetized by their communities

But AI is missing the critical infra layer needed  
to make this happen

# Building the Foundations for a Collaborative AI Economy

- To realize the promise of AI x Web3, we need more than theoretical alignment, we need real, usable infrastructure
- That means protocols for attribution, platforms for development, and systems that can scale across chains

# Building the Foundations for a Collaborative AI Economy

- To realize the promise of AI x Web3, we need more than theoretical alignment, we need real, usable infrastructure
- That means protocols for attribution, platforms for development, and systems that can scale across chains
- Sahara AI was designed to meet that need: not as a single app, but as an ecosystem foundation where open participation, verifiable contribution, and fair monetization are built in

# The Problem we are solving

Imagine you spent months curating a high-quality dataset or training an AI model. Today, once you share that dataset or model, you lose control over how it's used and receive no compensation when someone builds a billion-dollar AI application on top of your work.

# Economic Rails for Decentralized AI

Three Main Benefits:

# Economic Rails for Decentralized AI

Three Main Benefits:

01.



Fair Revenue:

Every contributor gets  
paid automatically  
based on their impact

# Economic Rails for Decentralized AI

Three Main Benefits:

01.



Fair Revenue:

Every contributor gets paid automatically based on their impact

02.



Verifiable Attribution:

Immutable on-chain proof of who built what

# Economic Rails for Decentralized AI

Three Main Benefits:

01.



Fair Revenue:

Every contributor gets paid automatically based on their impact

02.



Verifiable Attribution:

Immutable on-chain proof of who built what

03.



Composable Licensing:

Programmable usage rights that travel with the asset

# Why This Matters

Persona	Pain Points	Solution
Enterprises & Data Providers	Lose control and royalties once data is shared	Tokenize → License → Earn residuals recursively
AI Developers & Fine-tuners	Hard to monetize open-source contributions	Contribution NFTs + Staking Pool yield
Web3 Builders	No AI-native economic primitives	Plug-and-play NFTs, Vaults, Cashiers

# Real-World Flows: Registering and Monetizing a Dataset

# Real-World Flows: Registering and Monetizing a Dataset

01.



Registration:

Sarah uploads metadata and receives an AssetOwnerNFT with a Token-Bound Account.

# Real-World Flows: Registering and Monetizing a Dataset

01.



Registration:

Sarah uploads metadata and receives an AssetOwnerNFT with a Token-Bound Account.

02.



Setting Terms:

Sarah sets \$100 research and \$10K commercial licenses via TBA. Royalties: 60% hers, 20% team, 20% hospital.

# Real-World Flows: Registering and Monetizing a Dataset

01.



Registration:

Sarah uploads metadata and receives an AssetOwnerNFT with a Token-Bound Account.

02.



Setting Terms:

Sarah sets \$100 research and \$10K commercial licenses via TBA. Royalties: 60% hers, 20% team, 20% hospital.

03.



First Sale:

A pharma company buys a commercial license. The \$10K flows through the Cashier to the dataset's Vault.

# Real-World Flows: Registering and Monetizing a Dataset

01.



Registration:

Sarah uploads metadata and receives an AssetOwnerNFT with a Token-Bound Account.

02.



Setting Terms:

Sarah sets \$100 research and \$10K commercial licenses via TBA. Royalties: 60% hers, 20% team, 20% hospital.

03.



First Sale:

A pharma company buys a commercial license. The \$10K flows through the Cashier to the dataset's Vault.

04.



Automatic Distribution:

The Vault splits: \$6K to Sarah, \$2K to contributors, \$2K to the hospital.

# Real-World Flows: Registering and Monetizing a Dataset

01.



Registration:

Sarah uploads metadata and receives an AssetOwnerNFT with a Token-Bound Account.

02.



Setting Terms:

Sarah sets \$100 research and \$10K commercial licenses via TBA. Royalties: 60% hers, 20% team, 20% hospital.

03.



First Sale:

A pharma company buys a commercial license. The \$10K flows through the Cashier to the dataset's Vault.

04.



Automatic Distribution:

The Vault splits: \$6K to Sarah, \$2K to contributors, \$2K to the hospital.

05.



Downstream Revenue:

A model trained on Sarah's data earns revenue. She automatically gets 10% of its earnings, forever.

# Real-World Flows: Collaborative Model Development

# Real-World Flows: Collaborative Model Development

01.



Base Model  
Registration:

Team lead Alex registers the base model, receiving the AssetOwnerNFT.

# Real-World Flows: Collaborative Model Development

01.



Base Model  
Registration:

Team lead Alex registers the base model, receiving the AssetOwnerNFT.

02.



Contributions Flow In:

- Bob fine-tunes the model for medical applications
- Carol adds safety filters
- Diana adds data, each minted as a validated ContributionNFT.

# Real-World Flows: Collaborative Model Development

01.



Base Model  
Registration:

Team lead Alex registers the base model, receiving the AssetOwnerNFT.

02.



Contributions Flow In:

- Bob fine-tunes the model for medical applications
- Carol adds safety filters
- Diana adds data, each minted as a validated ContributionNFT.

03.



Revenue Sharing  
Setup:

The team sets splits via the Vault. Each ContributionNFT holder earns a share by contribution.

# Real-World Flows: Collaborative Model Development

01.



Base Model  
Registration:

Team lead Alex registers the base model, receiving the AssetOwnerNFT.

02.



Contributions Flow In:

- Bob fine-tunes the model for medical applications
- Carol adds safety filters
- Diana adds data, each minted as a validated ContributionNFT.

03.



Revenue Sharing  
Setup:

The team sets splits via the Vault. Each ContributionNFT holder earns a share by contribution.

04.



Ongoing Rewards:

As the model is used commercially, revenue automatically flows to all contributors according to their stakes.

# 4 Key Asset Ownership Questions

Sahara Protocol seeks to provide on-chain answers to four foundational asset questions:

Module Name	Definition
Who owns the asset?	Clear sovereignty and governance rights
Who contributed?	Contributor structure for attribution and future revenue participation
Who can use it?	On-chain usage rights and authorization proof
Is it legitimate?	On-chain verifiability rather than off-chain negotiation

# Token-Bound Accounts (TBAs): Giving Assets Their Own Wallets

# Token-Bound Accounts (TBAs):

## Giving Assets Their Own Wallets

01.



The asset can receive payments directly

02.



It can sign contracts and issue licenses

03.



It can pay out royalties to contributors

04.

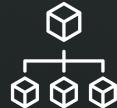


When you transfer ownership of the NFT, you transfer control of everything the asset owns

# The License System: Programmable Usage Rights



Commercial vs.  
non-commercial  
use



Whether  
derivatives are  
allowed



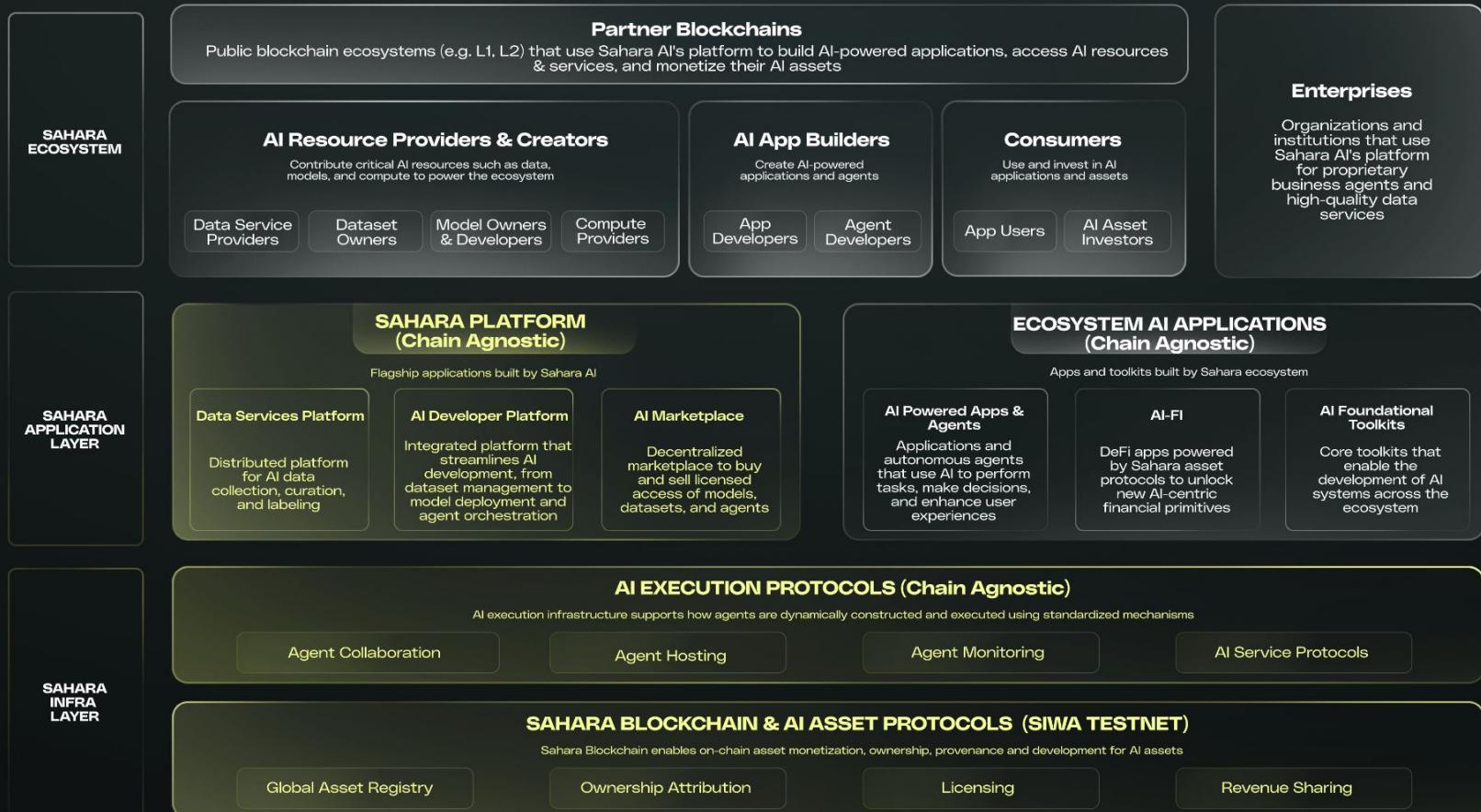
Time  
limitations



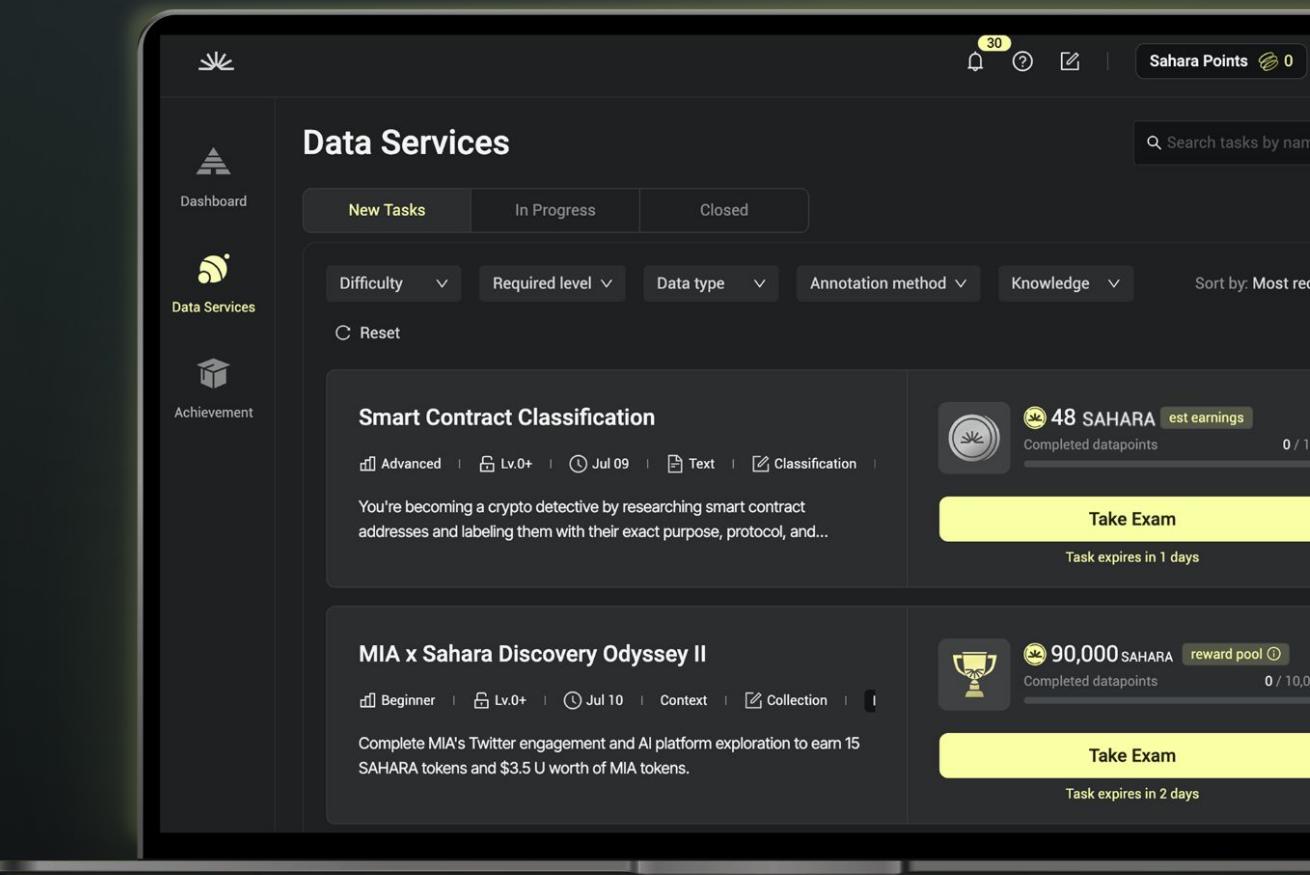
Geographic  
restrictions



Custom terms



# Data Services Platform



**Data Services**

New Tasks   In Progress   Closed

Difficulty   Required level   Data type   Annotation method   Knowledge

Sort by: Most recent

Reset

**Smart Contract Classification**

Advanced | Lv.0+ | Jul 09 | Text | Classification |

You're becoming a crypto detective by researching smart contract addresses and labeling them with their exact purpose, protocol, and...

48 SAHARA   est. earnings  
Completed datapoints   0 / 1

**MIA x Sahara Discovery Odyssey II**

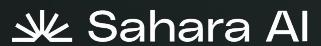
Beginner | Lv.0+ | Jul 10 | Context | Collection |

Complete MIA's Twitter engagement and AI platform exploration to earn 15 SAHARA tokens and \$3.5 U worth of MIA tokens.

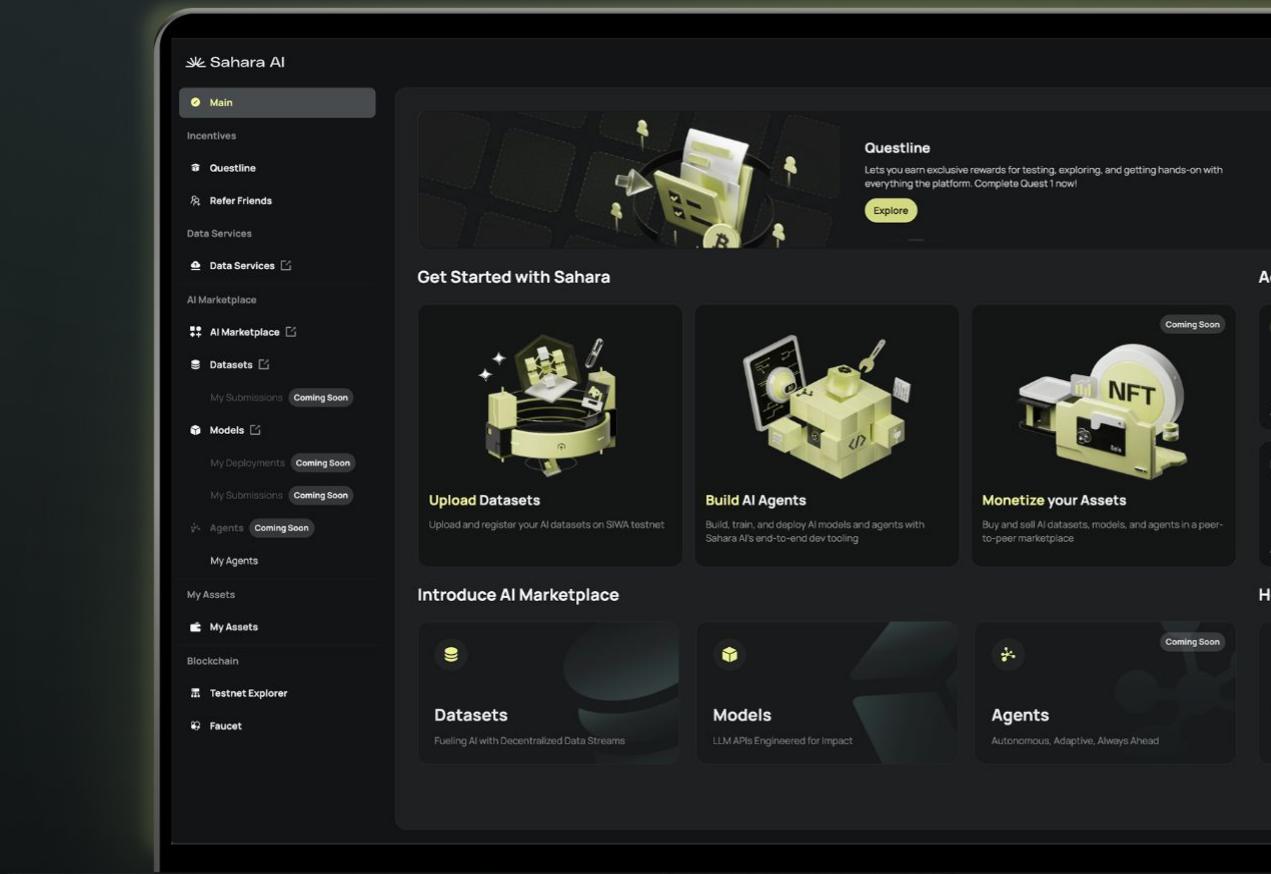
90,000 SAHARA   reward pool 0 / 10,000  
Completed datapoints   0 / 10,000

**Take Exam**  
Task expires in 1 days

**Take Exam**  
Task expires in 2 days



# AI Developer Platform



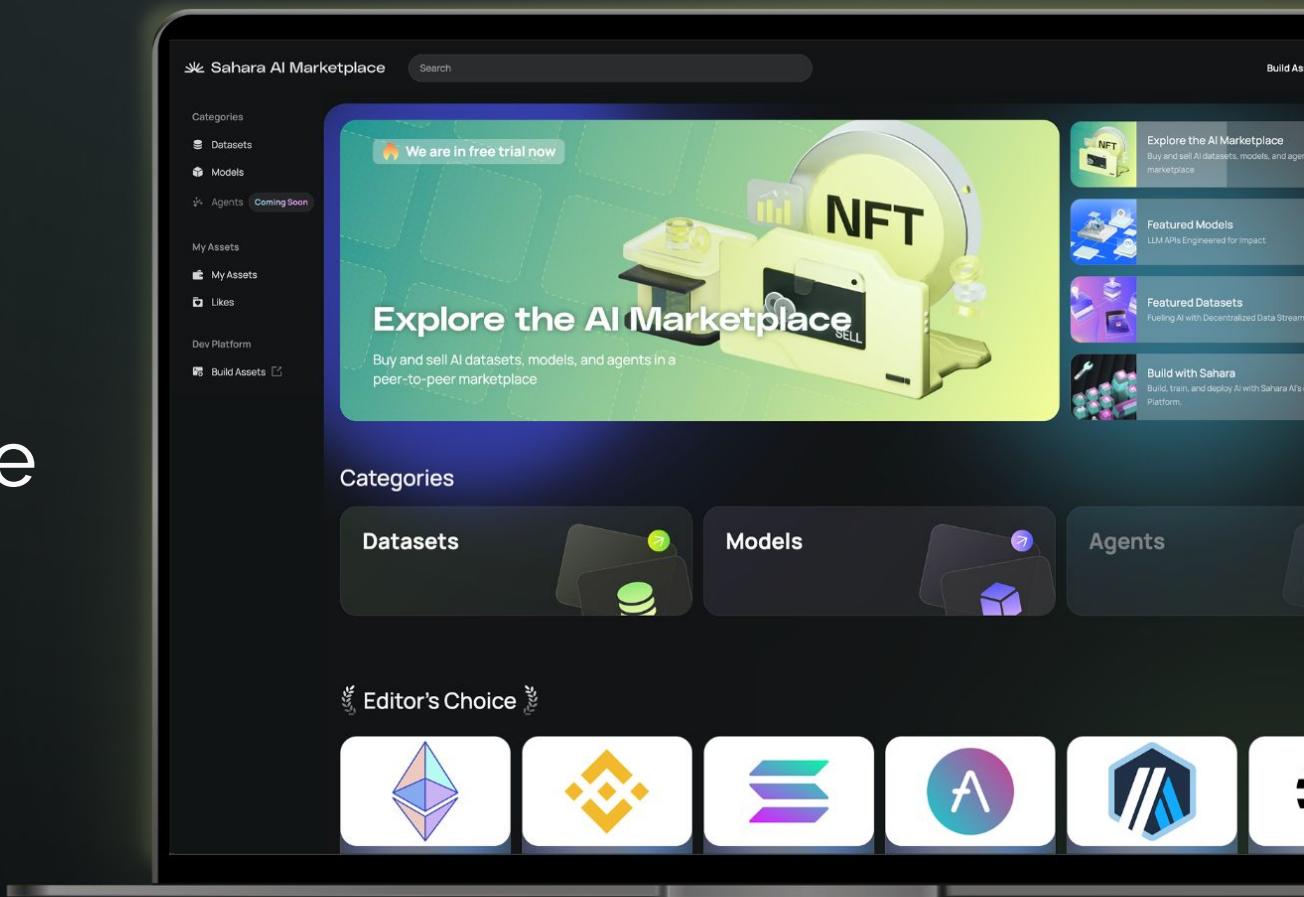
The image shows a screenshot of the Sahara AI platform interface. The left side features a dark sidebar with the following navigation items:

- Main (selected)
- Incentives
- Questline (Coming Soon)
- Refer Friends
- Data Services
  - Data Services (Coming Soon)
- AI Marketplace
- AI Marketplace (Coming Soon)
- Datasets (Coming Soon)
- My Submissions (Coming Soon)
- Models (Coming Soon)
  - My Deployments (Coming Soon)
  - My Submissions (Coming Soon)
- Agents (Coming Soon)
  - My Agents
- My Assets
- My Assets
- Blockchain
- Testnet Explorer
- Faucet

The main content area has a dark background with several sections:

- Questline**: Lets you earn exclusive rewards for testing, exploring, and getting hands-on with everything the platform. Complete Quest 1 now! [Explore](#)
- Get Started with Sahara**:
  - Upload Datasets**: Upload and register your AI datasets on SIWA testnet.
  - Build AI Agents**: Build, train, and deploy AI models and agents with Sahara AI's end-to-end dev tooling.
  - Monetize your Assets**: Buy and sell AI datasets, models, and agents in a peer-to-peer marketplace. (Coming Soon)
- Introduce AI Marketplace**:
  - Datasets**: Fueling AI with Decentralized Data Streams.
  - Models**: LLM APIs Engineered for Impact.
  - Agents**: Autonomous, Adaptive, Always Ahead. (Coming Soon)

# AI Marketplace



The image shows a screenshot of the Sahara AI Marketplace interface. The top navigation bar includes the logo, a search bar, and a "Build As" button. The left sidebar features a "Categories" section with "Datasets", "Models", and "Agents" (marked as "Coming Soon"), followed by "My Assets" (with "My Assets" and "Likes" sub-options), and a "Dev Platform" section with "Build Assets". The main content area has a large central banner with the text "We are in free trial now", "Explore the AI Marketplace", and "Buy and sell AI datasets, models, and agents in a peer-to-peer marketplace". The banner features a 3D rendering of a smartphone displaying an "NFT" logo, surrounded by various AI-related icons. Below the banner is a "Categories" section with three main items: "Datasets" (represented by a stack of three cylinders), "Models" (represented by a stack of three cubes), and "Agents" (represented by a stack of three hexagons). Further down is a "Editor's Choice" section featuring five icons: a purple and blue diamond shape, a yellow diamond shape, a purple and blue horizontal bar, a purple circle with a white "A" inside, and a blue hexagon with white diagonal lines. On the right side of the interface, there are three vertical columns of cards: "Explore the AI Marketplace" (with a green "NFT" icon), "Featured Models" (with a blue icon of a smartphone and laptop), "Featured Datasets" (with a purple icon of a smartphone and laptop), and "Build with Sahara" (with a blue icon of a smartphone and laptop).

Sahara AI Marketplace

Search

Build As

Categories

Datasets

Models

Agents Coming Soon

My Assets

My Assets Likes

Dev Platform

Build Assets

We are in free trial now

Explore the AI Marketplace

Buy and sell AI datasets, models, and agents in a peer-to-peer marketplace

NFT

Categories

Datasets

Models

Agents

Editor's Choice

NFT

Explore the AI Marketplace

Buy and sell AI datasets, models, and agents in a peer-to-peer marketplace

Featured Models

Featured Datasets

Build with Sahara

We must rewrite the relationship between  
AI companies, users, and contributors, to create a  
more equitable and transparent ecosystem for all.

An AI-Driven Future Is Inevitable.

An AI-Driven Future Is Inevitable.  
An Open and Equitable One Isn't.

# Join the Collaborative AI Economy